

Kriging on highly skewed data for DTPA-extractable soil Zn with auxiliary information for pH and organic carbon[☆]

J. Wu¹, W.A. Norvell^{*}, R.M. Welch

U.S. Plant, Soil and Nutrition Laboratory, USDA-ARS and Cornell Univ., Tower Road, Ithaca, NY 14853, USA

Received 4 March 2004; received in revised form 28 October 2005; accepted 9 November 2005

Available online 4 January 2006

Abstract

Knowledge of the distribution of crop-available trace elements in soils is limited by the sparseness of georeferenced data and the inherent variability of the more-labile forms of these elements. Cokriging with auxiliary variables can sometimes improve estimates for a less densely sampled primary variable, while skewed data can often be made more suitable for geostatistical modeling by appropriate transformation. Benefits from data transformation and cokriging in predicting Zn(DTPA) (an estimate of plant-available Zn, extracted from soil by the chelating agent diethylenetriaminepentaacetic acid) were assessed using a georeferenced set of data from northern North Dakota. Soil organic carbon (OC) and pH were used as auxiliary variables for cokriging. Data for Zn(DTPA), OC and pH were available for 587 locations. The statistical distribution of the data for Zn(DTPA) was highly skewed (approximately log-normal). Three methods of data transformation (computation of logarithms, conversion to standardized rank order and assignment of normal scores) were carried out prior to kriging or cokriging to reduce skewness. For comparisons of predictive success, the Zn(DTPA) data were partitioned into a predictor set of 293 sites and a testing set of 294 sites, according to a stratified randomized approach. Data for Zn(DTPA) in the testing set were reserved for testing estimates based on the predictor set. Cokriging on Zn(DTPA), using OC or pH as auxiliary variables, was consistently more effective than kriging on Zn(DTPA) alone. Cokriging with OC and pH together provided additional benefit. Data transformation generally improved kriged estimates, especially for low concentrations of Zn(DTPA) (e.g., $<0.5 \text{ mg kg}^{-1}$), which are important because they are indicative of soils containing inadequate Zn for optimal crop growth. Differences among normal score cokriging, log-normal cokriging and rank-ordered cokriging were relatively small.

Published by Elsevier B.V.

Keywords: Skewed distribution; Transformation; Zinc availability; Ordinary kriging; Log-normal; Rank order; Normal score; Cokriging; Auxiliary variables

Abbreviations: OK, ordinary kriging; OCK, ordinary cokriging; S.D., standard deviation; Zn(DTPA), DTPA-extractable (available) Zn in soils; OC, organic carbon; LG, logarithmic transform; RK, rank order transformation; NS, normal score transformation; cdf, cumulative distribution function; ME, mean error; RMSE, root mean square error.

[☆] Mention of proprietary product or vendors does not imply approval or recommendation by the U.S. Department of Agriculture.

^{*} Corresponding author. Fax: +1 607 255 1132.

E-mail address: WAN1@cornell.edu (W.A. Norvell).

¹ Current address: Department of Resources and Environment, Zhejiang Univ., Hangzhou, 310029, China.

1. Introduction

Knowledge of the availability to plants of trace elements in soils is important for maintaining or improving crop production and food quality. Several of the trace elements are micronutrients, required in small amounts by plants or animals for normal nutrition and health, and yet high concentrations of these or other trace elements can be toxic (Welch, 1995). Zinc is an example of an essential micronutrient that is often present in soils at levels that are inadequate for optimal nutrition of crops or for human consumers of crop-derived foods. Available Zn may be inadequate for optimal crop growth in as many as 50% of soils worldwide (Sillanpää, 1990), as well as in the diets of most humans, particularly in developing countries (Gibson, 1994).

Our understanding of the distribution of trace elements in soils is usually limited by a low geographic density of reliable data, as well as by high variability, skewed statistical distributions and unknown or poorly understood spatial dependencies. In contrast to the limited availability of information on trace elements, data related to major soil characteristics, such as organic carbon (OC), pH, cation-exchange capacity (CEC) or texture are often much more readily available. The spatial relationships between trace element concentrations and major soil characteristics can sometimes be used to improve predictions of the former, as we have shown for the prediction of total soil Cu in soils of northern North Dakota by cokriging with CEC as an auxiliary variable (Wu et al., 2003).

The distribution of crop-available Zn in soils of North Dakota is of interest because a substantial number of soils contain levels that are considered low or marginal for crop production based on crop responses and on soil testing extractions using the chelating agent DTPA (diethylenetriaminepentaacetic acid) (Franzen, 1999). The common occurrence of soils with low concentrations of DTPA-extractable Zn [Zn(DTPA)] in North Dakota is suggested also by our own results for almost 600 soils from the northernmost tier of 18 counties (Norvell et al., unpublished, 2002). These results suggest that maps of Zn(DTPA) in this important agricultural region would be useful to farmers and agricultural extension personnel in identifying soils that might benefit from Zn fertilization. Preliminary analyses of the data from this survey showed that Zn(DTPA) was related to OC and pH, suggesting that these characteristics might serve effectively as auxiliary variables to improve geostatistical estimates of the distribution of available Zn in soils. These relationships

could provide an important benefit in cokriging analyses and map preparation because data for OC and pH are available from other sources (e.g., Holmgren et al., 1993; National Soil Survey Center, 2002).

Geostatistical inferences using kriging techniques are more efficient when data for variables are distributed normally. However, data for concentrations of elements in soils or geologic materials are often skewed or highly skewed (e.g., Journel, 1980; McBratney et al., 1982; Juang et al., 2001; Webster and Oliver, 2001). Our data for Zn(DTPA) in soils of North Dakota are no exception, as we will show below. Difficulties caused by highly skewed distributions can often be alleviated by appropriate transformation of the data. The most common is the logarithmic transform (Journel, 1980; Saito and Goovaerts, 2000), which is best suited to log-normal data. A second approach is to use a standardized rank order transformation prior to kriging, a simple method that is well suited to integrating many diverse types of data (Journel and Deutsch, 1997). A third approach is normal score transformation (Goovaerts, 1997; Deutsch and Journel, 1998), a procedure that transforms any distribution into a normalized distribution. All three of the above transformations appeared potentially suited to improving the analysis of our soil data for Zn(DTPA) and all were utilized in this study.

Other approaches for handling skewed data include indicator kriging, multiple indicator kriging (Deutsch and Journel, 1998; Saito and Goovaerts, 2000; Cattle et al., 2002) and disjunctive kriging (Van Meirvenne and Goovaerts, 2001; Webster and Oliver, 2001). These methods were considered, but not pursued in our work involving cokriging with auxiliary variables. The principal reason for excluding these methods was practical rather than conceptual in that the necessary calculations using both primary and secondary variables are sufficiently burdensome so as to be impractical with currently available geostatistical programs. An additional disadvantage to indicator (co)kriging is that much information is lost in the categorization of originally continuous data into a single indicator. And, in regard to disjunctive kriging, Deutsch and Journel (1998, p. 17) have suggested that this method may be replaced by the more robust approach of kriging following normal-score transformation of data, which is one of the methods that we have used below.

The objectives of our study were to utilize available data from soils of northern North Dakota to: (i) compare four kriging methods [ordinary (co)kriging, log-normal ordinary (co)kriging, rank order ordinary (co)kriging and normal score ordinary (co)kriging] in their predictions for Zn(DTPA), a measure of a crop-avail-

able trace element with highly skewed data; and (ii) determine if the kriged estimates of concentrations of Zn(DTPA) by these methods can be improved by incorporating information on the auxiliary soil characteristics, OC and pH, when data for Zn(DTPA) are not available.

2. Data set and methods

2.1. Data set

Data for Zn(DTPA), OC and pH were available from our concurrent study of the geographical distribution of micronutrients and trace elements in the 18 counties of northern North Dakota (Norvell et al., unpublished data, 2002). These results were obtained from soils suited to crop production, sampled according to a stratified randomized design (Petersen and Calvin, 1986) in combination with nested sampling to characterize short-range variations. A total of 587 sites with complete data

for Zn(DTPA), OC and pH were selected for our current objectives. Concentrations of Zn(DTPA) were measured according to Lindsay and Norvell (1978). The analytical methods for determining pH and OC were described by Norvell et al. (2000).

The Zn(DTPA) data for the 587 sites were partitioned into two subsets, a predictor subset and a testing subset. To provide stratification, all sites were first sorted by their site-numbers, which had been assigned generally according to their geographic locations from west to east and from north to south within a county. Then, for every two sites, one was randomly selected as a predictor site and the other as testing site. This process partitioned the full set of 587 sites into a subset of 293 for predictions and a subset of 294 for testing. The geographic distributions of these sites are shown in Fig. 1. Summary statistics are given in Table 1. Strongly positive skewness, approximating log-normality, is shown by this Zn(DTPA) data and by the frequency histogram in Fig. 2.



Fig. 1. Locations of 587 sites with measured concentrations of Zn(DTPA), pH and organic carbon in 18 counties of northern North Dakota, USA.

Table 1

Summary statistics of the attributes for the full data set, the predictor subset and the testing subset used in the study

	Zn(DTPA) ^a	Log Zn(DTPA) ^a	OC ^a	pH	Zn (DTPA)	Log Zn(DTPA)	OC	pH	Zn(DTPA)	Log Zn(DTPA)	OC	pH
	Full set (n=587)				Predictor set (n=293)				Testing set (n=294)			
Mean	0.99	−0.123	23.5	7.38	0.96	−0.123	23.6	7.39	1.01	−0.122	23.4	7.37
Median	0.73	−0.137	22.8	7.63	0.73	−0.137	23.1	7.67	0.72	−0.143	22.7	7.61
S.D. ^b	0.87	0.311	8.9	0.78	0.76	0.301	8.6	0.78	0.96	0.322	9.19	0.79
Min ^b	0.10	−1.000	2.7	4.97	0.10	−1.000	4.0	4.97	0.13	−0.886	2.7	5.23
Max ^b	9.15	0.961	68.8	8.96	5.38	0.731	68.8	8.83	9.15	0.961	59.0	8.96
Skew ^b	3.31	0.18	0.57	−0.72	2.39	0.10	0.70	−0.74	3.65	0.26	0.47	−0.71
Kurt ^b	18.94	0.05	1.08	−0.39	8.32	−0.09	2.15	−0.34	21.41	0.16	0.26	−0.42

^a Zn(DTPA), DTPA-extractable Zn (mg kg^{−1}); log Zn(DTPA), logarithm of Zn(DTPA); OC, organic C (g kg^{−1}).^b S.D., standard deviation; min, minimum; max, maximum; skew, skewness; kurt, kurtosis.

2.2. Kriging methods

Before kriging, the spatial variability of original or transformed variables was modeled with the aid of (cross)semivariograms, which are graphs of (cross)semivariances, $\gamma_{Z_i Z_j}(h)$, as a function of the separation or lag distance, h , between all possible pairs of sample locations. The (cross)semivariances were computed using:

$$\gamma_{Z_i Z_j}(h) = \frac{1}{2n(h)} \sum_{\alpha=1}^{n(h)} [Z_i(u_\alpha) - Z_i(u_\alpha + h)] \times [Z_j(u_\alpha) - Z_j(u_\alpha + h)] \quad (1)$$

where Z_i and Z_j are variables; $Z(u)$ and $Z(u+h)$ are two observations of a variable separated by the lag distance, h ; and $n(h)$ is the number of pairs of observations at this lag distance.

We found that semivariograms for the three variables could be fitted successfully with either spherical or exponential mathematical models. However, because of the need to model all variograms simultaneously in cokriging, we found the spherical model was more

effective in fitting data with a single set of parameters, thereby meeting the requirements of semi-positive definiteness (Goovaerts, 1999; Goulard and Voltz, 1992). Thus, linear combinations of an isotropic spherical model were used to describe the spatial variability and fit the variograms:

$$\gamma(h) = 0 \quad \text{for } h = 0$$

$$\gamma(h) = C_0 + C_1 \left[1.5(h/A_1) - 0.5(h/A_1)^3 \right] + C_2 \left[1.5(h/A_2) - 0.5(h/A_2)^3 \right] \quad \text{for } h \leq A_1$$

$$\gamma(h) = C_0 + C_1 + C_2 \left[1.5(h/A_2) - 0.5(h/A_2)^3 \right] \quad \text{for } h \leq A_2$$

$$\gamma(h) = C_0 + C_1 + C_2 \quad \text{for } h > A_2 \quad (2)$$

where C_0 is nugget variance; C_1 is the first structural variance at the (first) range distance, A_1 ; and C_2 is the second structural variance at A_2 , assuming $A_2 > A_1 > 0$.

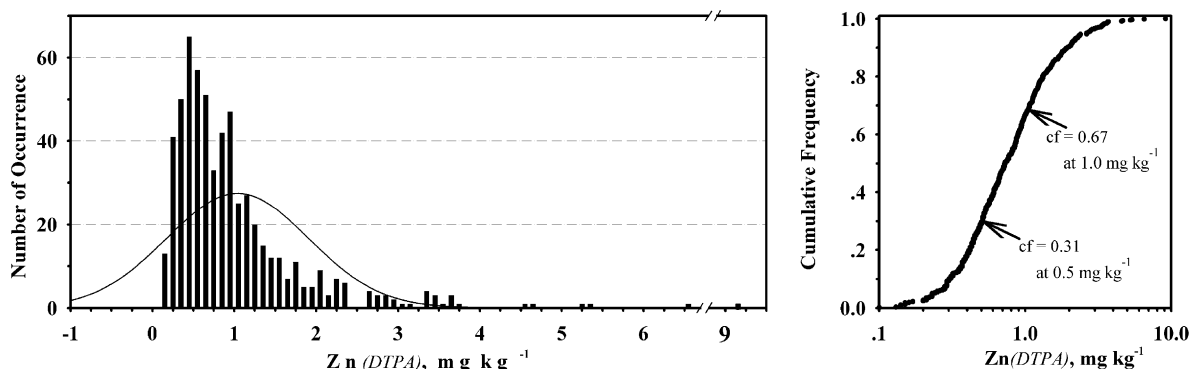


Fig. 2. Frequency histogram of Zn(DTPA) data (left panel) and its cumulative frequency distribution (right panel). The bell-shaped line on the histogram shows a normal distribution with the same mean (0.99 mg kg^{−1}) and standard deviation (0.87 mg kg^{−1}) as the measured population.

The ordinary kriging (OK) estimator is expressed as:

$$Z_{OK}^*(u) = \sum_{\alpha=1}^{n(u)} \lambda_{OK}(u_{\alpha}) Z(u_{\alpha}) \quad (3)$$

where $Z_{OK}^*(u)$ is the estimated value of Z at location u ; $\lambda_{OK}(u_{\alpha})$ corresponds to the weight associated with the measured value of Z at location u_{α} . The weights are determined so that the estimated error variance is minimized. The λ_{OK} are forced to $\sum_{\alpha=1}^{n(u)} \lambda_{OK}(u_{\alpha}) = 1$, in which $n(u)$ is the number of measured values used in estimation in neighborhoods of u .

The ordinary cokriging (OCK) estimator for the primary variable with multiple secondary variables is:

$$Z_{OCK}^{(1)*}(u) = \sum_{\alpha_1=1}^{n_1(u)} \lambda_{OCK}(u_{\alpha_1}) Z_1(u_{\alpha_1}) + \sum_{i=2}^{N_v} \sum_{\alpha_i=1}^{n_i(u)} \lambda_{OCK}(u) Z_i(u_{\alpha_i}) \quad (4)$$

where $Z_{OCK}^{(1)*}(u)$ is the value of the primary variable Z_1 to be estimated at the location u and $\lambda_{OCK}(u_{\alpha_i})$ corresponds to the weight associated with the measured values of Z_i at the location. As described above, the weights are selected to minimize the estimated error variance u_{α_i} . Weights for the primary variable $\lambda_{OCK}(u_{\alpha_i})$ are forced to $\sum_{\alpha_1=1}^{n_1(u)} \lambda_{OCK}(u_{\alpha_1}) = 1$, while for the secondary variables $\sum_{\alpha_i=1}^{n_i(u)} \lambda_{OCK}(u_{\alpha_i}) = 0$; N_v is the number of variables; and $n_i(u)$ is the number of neighboring values, Z_i , used in estimating u .

In this study, all kriging inferences were made using GSLIB (Deutsch and Journel, 1998). A search radius of 50 km was used, with a minimum number of 12 points for the primary variable and a maximum number of 24 points for primary or secondary variables in all interpolations. More detailed descriptions of the geostatistical methods can be found in text books such as Goovaerts (1997) and Webster and Oliver (2001).

2.3. Kriging on transformed data

Transformation of data may be desirable before kriging to normalize the data distribution, suppress outliers and improve data stationarity. Appropriate transformation may also make spatial relations more evident and provide more stable variograms. Because of the skewness of our Zn(DTPA) data, we selected logarithmic, rank order and normal score transformations as means to improve the predictions of Zn(DTPA) by kriging or cokriging. Because we wished to interpret results for Zn(DTPA) in their original concentration

units, back-transformation of our (co)kriged results was necessary.

Log-normal ordinary kriging or cokriging (OK_{LG} or OCK_{LG}) is performed on log-transformed data, i.e.,

$$y(u_{\alpha}) = \log z(u_{\alpha}) \quad (5)$$

where $z(u_{\alpha})$ is the measured value at location u_{α} . Back-transformation of each (co)kriging result was carried out by exponentiation to reverse Eq. (5), providing a prediction for Zn(DTPA) expressed in original concentration units.

Rank order ordinary kriging or cokriging (OK_{RK} or OCK_{RK}) is performed on standardized rank order transformed data. For a variable Z with a cumulative distribution function (cdf), $F(z)$, the transformed data have a uniform distribution on the interval from zero to one. For a set of n samples, the empirical cumulative distribution function is used to estimate $F(z)$. In practice, the transformation and back-transformation are carried out as follows (Journel and Deutsch, 1997; Juang et al., 2001):

1. Arrange the n sample in ascending order:

$$z^{(1)} \leq \dots \leq z^{(r)} \leq \dots \leq z^{(n)} \quad (6)$$

where the superscript r is the rank of datum $z^{(r)}$ among all n data, $z^{(r)}$ is called the r th order statistic.

2. Calculate the standardized rank $y^{(r)}$ of the sample

$$y^{(r)} = \frac{r}{n} \quad (7)$$

The value of $y^{(r)}$ is between $1/n$ and 1.

3. Kriging is carried out on the ranks. Estimated ranks, $y^*(u)$, are back-transformed into the original units for variable Z :

$$z^*(u) = F^{-1}(y^*(u)) \quad (8)$$

Most estimated values for $y^*(u)$ usually fall between two adjacent standardized ranks, say r/n and $(r+1)/n$. Under the circumstances, the corresponding estimates in the original concentration space $z^*(u)$ will be between $z^{(r)}$ and $z^{(r+1)}$. Thus, the value of $z^*(u)$ is assigned to the mid-point between $z^{(r)}$ and $z^{(r+1)}$ (Juang et al., 2001):

$$z^*(u) = 0.5 [z^{(r)} + z^{(r+1)}] \quad (9)$$

If $y^*(u)$ happens to be r/n , then

$$z^*(u) = z^{(r)} \quad (10)$$

On occasion, a value for $y^*(u)$ estimated by kriging may fall outside the acceptable range between the minimum of $1/n$ and the maximum of 1. In this case,

we re-assigned any estimate $<1/n$ to equal $1/n$ and any estimate >1 to equal 1, prior to back-transformation.

Normal score ordinary kriging or cokriging (OK_{NS} or OCK_{NS}) is performed on the normal score transformed data. The normal score transform is a graphical transform similar to a correspondence table between equal p quantiles z_p and y_p of the Z cdf $F(z)$ and the standard Gaussian cdf $G(y)$. The normal score transform is carried out as follows (Deutsch and Journel, 1998; Saito and Goovaerts, 2000):

1. The n sample data are ranked in ascending order similar to Eq. (6):

$$z^{(1)} \leq \dots \leq z^{(k)} \leq \dots \leq z^{(n)} \quad (11)$$

where the superscript k is the rank of datum $z^{(k)}$ among all n data;

2. The sample cumulative frequency of the datum $z^{(k)}$ is then computed as:

$$p^k = k/n \quad (12)$$

3. The normal score transform of the $z^{(k)}$ datum is matched to the p^k quantile of the standard normal cdf:

$$y^{(k)} = G^{-1}\left\{F\left[z^{(k)}\right]\right\} = G^{-1}(p^k) \quad (13)$$

4. Kriging is performed on the transformed data. Estimates of the standard normal deviate, $y^*(u)$, are back-transformed to original units:

$$z^*(u) = F^{-1}(G(y^*(u))) \quad (14)$$

where $F(z)$ is the cdf of the original data.

The back-transformation of (co)kriged results by the simple reversal of the transformation procedures described above, converts values representing kriged means (in transformed units), into median predictions (in the original units of concentration). Correcting for the reduction in skewness introduced by back-transformation into median values is not straightforward. Empirical corrections can be devised, but their suitability is uncertain (Saito and Goovaerts, 2000). Consequently, we chose back-transformation based on straight-forward reversal of the three transformations, assuming that the benefits from carrying out (co)kriging with data that more closely met the assumption of stationarity would outweigh any bias introduced by utilizing predicted median concentrations of Zn(DTPA) to represent testing locations.

2.4. Evaluation of kriging methods

To evaluate the performance of the kriging methods, the data for Zn(DTPA) from the 294 sites in the testing

set were excluded from kriging and reserved for testing accuracy of predictions. Data for Zn(DTPA) at the 293 sites in the predictor data set, along with pH and OC data for all 587 sites, were then used in geostatistical inferences leading to kriged estimates of Zn(DTPA) at the sites reserved for testing. Descriptive statistics and scatter plots were used to compare true (measured) concentrations of Zn(DTPA) with the back-transformed predictions of the several kriging methods. In addition, we computed the mean error (ME), root mean square error (RMSE) and the coefficient of determination (r_p^2). The ME and RMSE have their standard meanings (Isaaks and Srivastava, 1989) and the definition of r_p^2 follows Juang and Lee (1998):

$$ME = \frac{1}{n} \sum_{i=1}^n [z(u_i) - z^*(u_i)] \quad (15)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n [z(u_i) - z^*(u_i)]^2} \quad (16)$$

$$r_p^2 = 1 - \frac{\sum_{i=1}^n [z^*(u_i) - z(u_i)]^2}{\sum_{i=1}^n [z(u_i) - \bar{z}(u)]^2} \quad (17)$$

where $z(u_i)$ is the measured value of Z at location u_i , $z^*(u_i)$ is the predicted value at the same location and $\bar{z}(u)$ is the mean of the measured values. The ME provides a measure of bias; the RMSE provides a measure of accuracy; and the coefficient of determination, r_p^2 , expresses the proportion of the original variance accounted for by the model predictions.

3. Results and analyses

Summary statistics show that Zn(DTPA), OC and pH all varied substantially among the 587 sites (Table 1, Fig. 2). The range in Zn(DTPA) was especially large, 0.10 to 9.15 mg kg⁻¹, while pH varied from 4.97 to 8.96, and OC ranged from 2.7 to 68.8 g kg⁻¹. The wide range in Zn(DTPA) included many low concentrations suggesting a low availability of Zn to crops. The 587 sites included 176 sites with Zn(DTPA) <0.5 mg kg⁻¹ and 215 sites with Zn(DTPA) between 0.5 and 1.0 mg kg⁻¹, indicating respectively soils with insufficient or marginal supplies of Zn for sensitive crops in North Dakota (Franzen, 1999).

The data for OC and pH had low skewness and kurtosis, but data for Zn(DTPA) were far from normally distributed (Fig. 2). For Zn(DTPA), the mean and median differed by more than a third, the S.D. exceeded

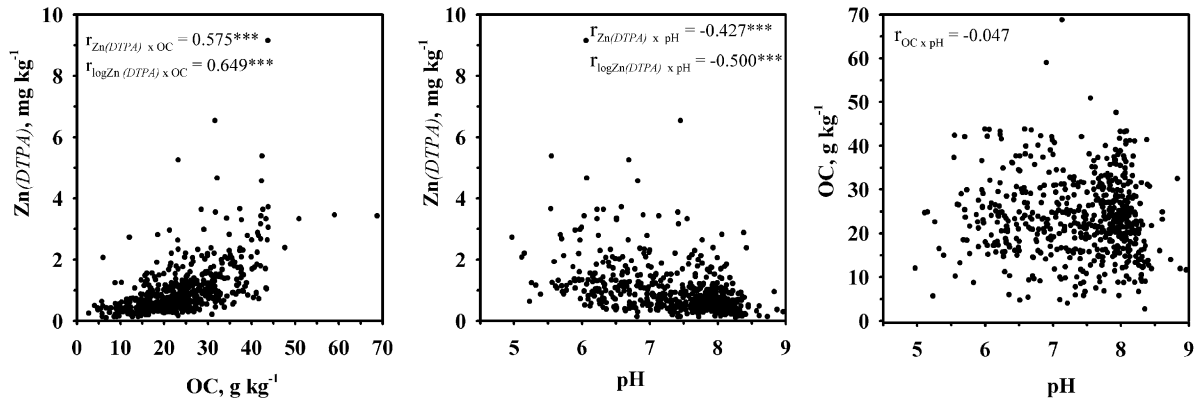


Fig. 3. Scattergrams among Zn(DTPA), organic carbon (OC) and pH in 587 soils in northern North Dakota.

the median and the data were strongly positively skewed with high kurtosis. Data for Zn(DTPA) were distributed more-nearly in a log-normal manner. A common log-transformation was successful in normalizing the data as shown by the reduction of skewness and kurtosis to 0.18 and 0.05, respectively (Table 1). The summary statistics for rank order and normal score transformations are not shown, because these values are

set by transformation to their standardized values (Journel and Deutsch, 1997; Deutsch and Journel, 1998).

Comparing the predictor and testing subsets shows that their summary statistics are quite close for log Zn(DTPA) and pH, and reasonably close for OC. For untransformed Zn(DTPA), however, the differences between subsets were greater. These differences between subsets are an accidental result of randomized

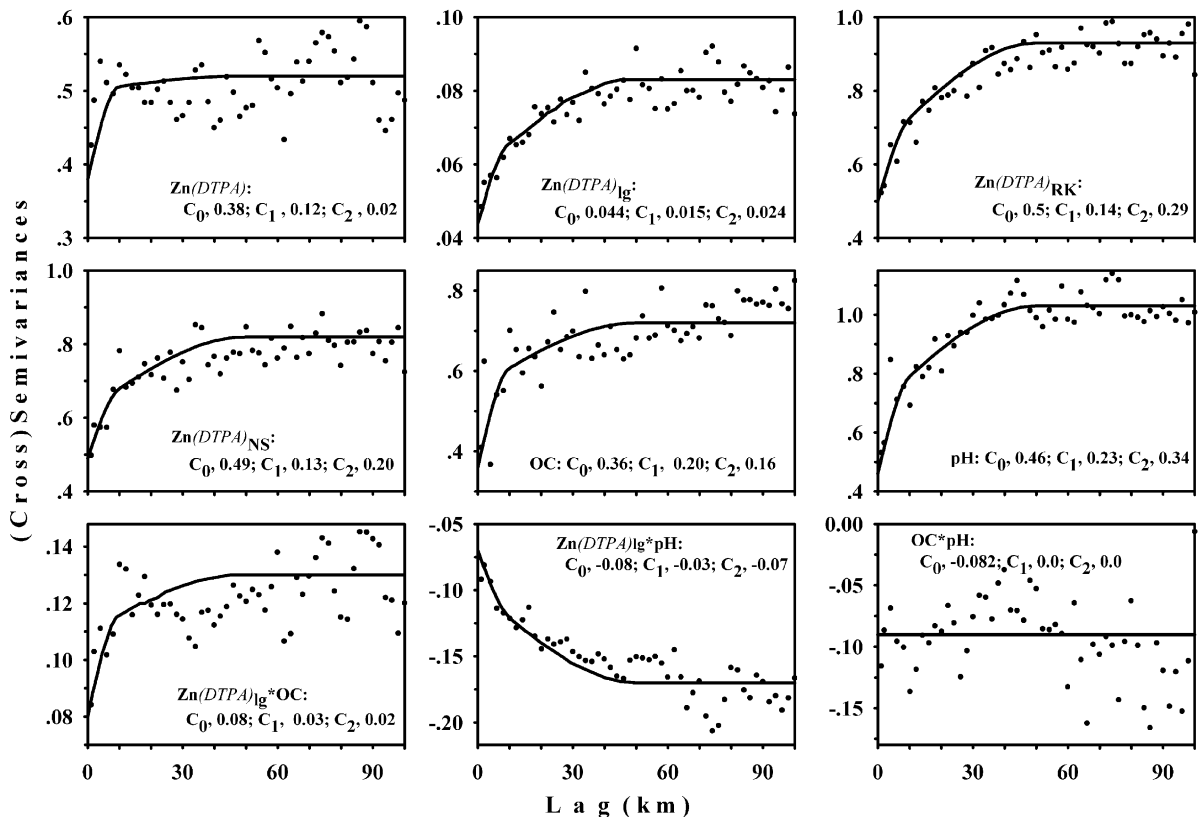


Fig. 4. Selected experimental omnidirectional (cross)-semivariograms, and fitted curves and parameters of the double spherical model (see Eq. (2) in the text). Range distances A_1 and A_2 are set for 10 km and 50 km correspondingly for all variogram models.

Table 2

Fitted parameters^a for the double spherical model (see Eq. (2)) for six experimental cross-semivariograms which are not shown in Fig. 4

Variables ^b	C_0	C_1	C_2
OC * Zn(DTPA)	0.18	0.07	0.05
OC * Zn(DTPA) _{RK}	0.24	0.12	0.04
OC * Zn(DTPA) _{NS}	0.28	0.10	0.03
pH * Zn(DTPA)	−0.20	−0.09	−0.08
pH * Zn(DTPA) _{RK}	−0.22	−0.12	−0.22
pH * Zn(DTPA) _{NS}	−0.22	−0.12	−0.22

^a C_0 , nugget variance; C_1 , structural variance at range distance A_1 ; C_2 , structural variance at range distance A_2 ; range distances A_1 and A_2 were set at 10 and 50 km, correspondingly, for all variogram models. Other fitted variance parameters are shown in Fig. 4.

^b OC, organic C; Zn(DTPA), DTPA-extractable Zn; Zn(DTPA)_{RK}, rank-order transformed Zn(DTPA); Zn(DTPA)_{NS}, normal score transformed Zn(DTPA).

partitioning of the data, with much of the difference arising from the less-than-equal distribution of extreme values that commonly occurs in randomized partitions of highly skewed populations. To the extent that these extreme values are more difficult to model, they are a challenge for each of the kriging methods compared.

The bivariable relations among the variables are shown in scattergrams (Fig. 3). The correlation between

OC and Zn(DTPA) or log Zn(DTPA) was moderately good, while the correlation between pH and either of these variables was weaker although still highly statistically significant. Virtually, no correlation between OC and pH was observed, indicating that these variables explained different portions of the variability in Zn(DTPA). A multiple-linear regression resulted in:

$$\text{Zn(DTPA)} = 2.98 - 0.44\text{pH} + 0.054\text{OC}$$

$$R^2 = 0.49*** \quad (18)$$

where OC is in g kg^{-1} and Zn in mg kg^{-1} . The benefit from incorporating both OC and pH are shown by the significant increase in R^2 . Similar or greater increases in R^2 were found when OC and pH were incorporated in multiple-linear regressions for any of the transformed concentrations of Zn(DTPA).

The experimental (cross)variograms show different degrees of spatial auto-continuity or cross-continuity of attributes (Fig. 4). The variogram for Zn(DTPA) is scattered and erratic at all lag distances. This effect is attributable primarily to the influence of extreme values in the skewed data set. Well-behaved variograms with clear spatial structure were obtained for Zn(DTPA) following log, rank order or normal score transforma-

Table 3

Summary statistics^a for true and estimated Zn(DTPA) concentrations in soils (mg kg^{-1}) at the 294 sites of the testing set as predicted by four kriging methods^b with or without use of auxiliary variables, OC and pH

	Zn(DTPA) only				Zn(DTPA) with OC				Zn(DTPA) with pH				Zn(DTPA) with OC and pH				
	True	OK	OK _{LG}	OK _{RK}	OK _{NS}	OCK	OCK _{LG}	OCK _{RK} ^c	OCK _{NS}	OCK	OCK _{LG}	OCK _{RK} ^c	OCK _{NS}	OCK	OCK _{LG}	OCK _{RK} ^c	OCK _{NS}
Mean	1.01	0.95	0.83	0.80	0.82	0.93	0.85	0.89	0.86	0.96	0.87	0.87	0.87	0.93	0.88	0.88	0.88
Med	0.72	0.84	0.74	0.74	0.72	0.83	0.69	0.71	0.68	0.82	0.71	0.72	0.71	0.83	0.70	0.71	0.68
S.D.	0.96	0.45	0.41	0.31	0.41	0.60	0.60	0.77	0.63	0.58	0.52	0.60	0.51	0.64	0.61	0.75	0.63
Min	0.13	0.41	0.32	0.42	0.35	−0.44	0.24	0.23	0.23	−0.03	0.24	0.28	0.25	−0.55	0.17	0.10	0.15
Max	9.15	3.10	2.87	2.24	2.81	3.42	4.51	5.38	5.25	3.30	3.37	5.38	3.30	3.54	5.04	5.38	5.25
<i>q</i> _{0.25}	0.45	0.68	0.57	0.61	0.57	0.52	0.50	0.56	0.49	0.59	0.55	0.60	0.54	0.47	0.50	0.54	0.49
<i>q</i> _{0.75}	1.17	1.04	0.91	0.87	0.89	1.17	0.96	0.90	0.95	1.21	1.03	0.93	1.00	1.27	0.99	0.93	0.99
cf _{0.5}	0.31	0.05	0.07	0.06	0.15	0.24	0.23	0.19	0.27	0.18	0.19	0.13	0.19	0.27	0.27	0.18	0.27

Predictive success, all 294 sites

ME	0.06	0.18	0.21	0.19	0.08	0.16	0.12	0.15	0.05	0.14	0.14	0.14	0.14	0.08	0.13	0.13	0.13
RMSE	0.79	0.82	0.85	0.82	0.67	0.63	0.62	0.63	0.77	0.76	0.76	0.76	0.76	0.66	0.62	0.67	0.61

Predictive success, 92 sites $< 0.5 \text{ mg kg}^{-1}$

ME	−0.40	−0.29	−0.31	−0.29	−0.16	−0.17	−0.21	−0.17	−0.21	−0.20	−0.24	−0.21	−0.21	−0.03	−0.11	−0.15	−0.12
RMSE	0.46	0.36	0.36	0.35	0.36	0.27	0.27	0.26	0.36	0.27	0.29	0.27	0.27	0.31	0.20	0.21	0.19

^a Med, median; S.D., standard deviation; min, minimum; max, maximum; ME, mean errors; RMSE, root mean square errors; $q_{0.25}$, the lower quartile (mg kg^{-1}); $q_{0.5}$, the middle quartile (median) (mg kg^{-1}); $q_{0.75}$, the upper quartile (mg kg^{-1}); cf_{0.5}, cumulative frequency of Zn(DTPA) concentrations at $< 0.5 \text{ mg kg}^{-1}$ soil.

^b OK, ordinary kriging; OK_{LG}, log-normal ordinary kriging; OK_{RK}, rank-order ordinary kriging; OK_{NS}, normal score ordinary kriging; Ock, ordinary cokriging; Ock_{LG}, log-normal ordinary cokriging; Ock_{RK}, rank-order ordinary cokriging; Ock_{NS}, normal score ordinary cokriging.

^c For Ock_{RK}, a few cokriging estimates fell above the upper boundary in standardized rank (i.e., values > 1.0). These overestimates have no corresponding back-transform within the defined range for Zn concentrations, and thus were assigned to the maximum value of 1.0 as is customary (see Methods). To a small, but unknown extent, this restriction of high values is likely to artificially inflate predictive success by Ock_{RK}.

tion. Each of the three transformations successfully reduced the distortions from skewness and extreme values. The cross-variogram between OC and pH shows little spatial co-variability, which echoes the lack of relationship shown in their scatterplot (Fig. 3) and near-zero correlation coefficient.

The double spherical model (Eq. (2)) was chosen to provide reasonable fit for all the variograms, using A_1 set at 10 km and A_2 at 50 km. The fitted nugget variances and structural variances are shown either on the variograms of Fig. 4 or in Table 2. With these parameters, each coregionalization matrix was positive semi-definite as required to conduct cokriging (Govaerts, 1999).

Estimates for Zn(DTPA) at the 294 locations of the testing set were obtained by the four kriging methods

using the fitted parameters in Fig. 4 and Table 2. Back-transformations for the estimates based on log-normal ordinary (co)kriging, rank order (co)kriging or normal-score (co)kriging were performed according to formulas presented in Section 2.2.

Summary statistics for Zn(DTPA) estimated by all kriging methods for the 294 testing sites are tabulated in Table 3. For comparison, this table includes also the statistics for the true values of Zn at these same sites, i.e., the values that were excluded from kriging and reserved by design to validate estimates based on the predictor set data. Detailed comparisons between true values and estimates by the kriging methods are shown in Fig. 5. Perfect agreement between true and predicted values would be reflected in having all pairs of points fall on the 1:1 line (dashed line).

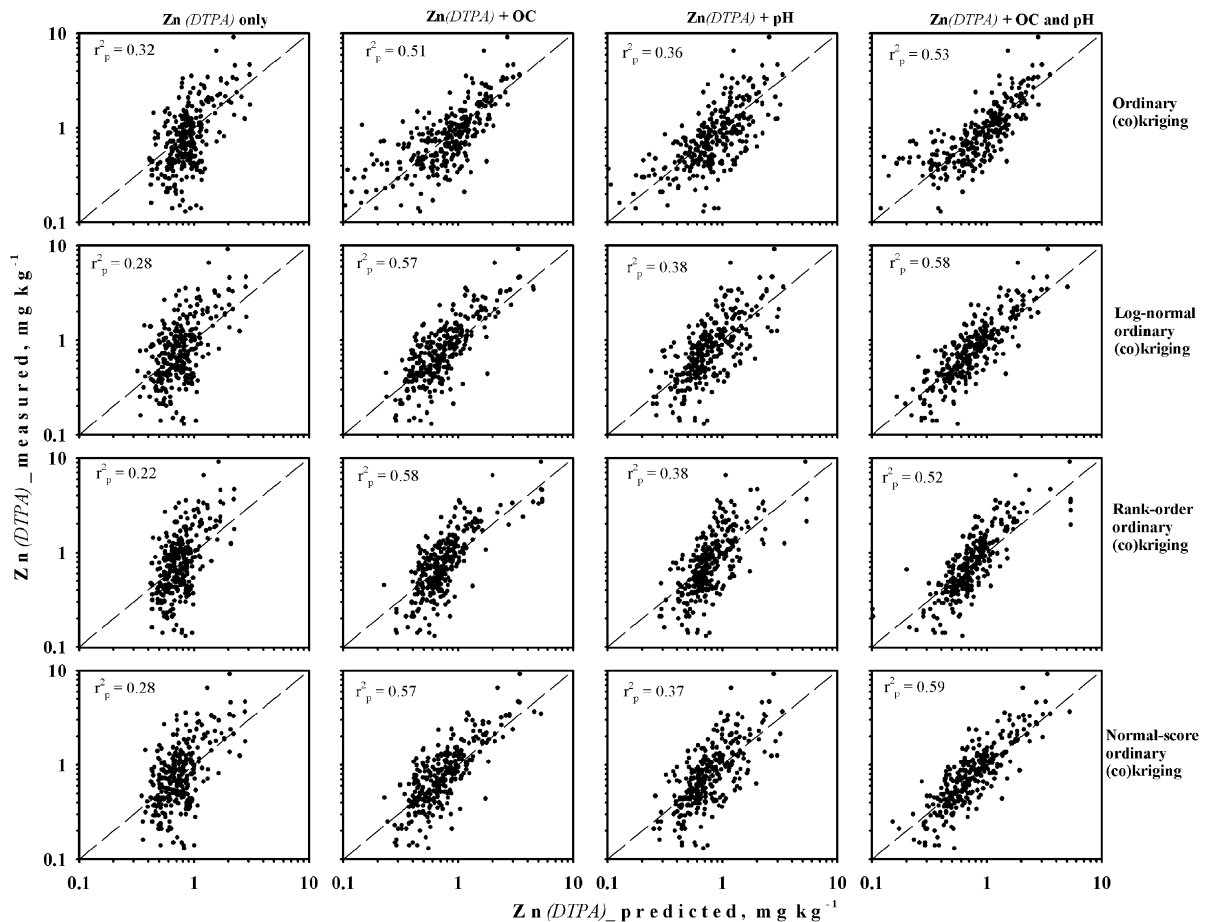


Fig. 5. Comparisons of measured (true) Zn(DTPA) concentrations in soils of the testing set ($n=294$) with values predicted by four methods of kriging or cokriging methods using four combinations of variables. From left to right, the first column shows results for kriging using Zn(DTPA) data only, while the plots in subsequent columns show cokriging results utilizing Zn(DTPA) and auxiliary data for organic carbon (OC), pH or both. From top to bottom, the four rows contain results from ordinary (co)kriging, log-normal ordinary (co)kriging, rank order ordinary (co)kriging and normal score ordinary (co)kriging, respectively. Both x-axis and y-axis use log scale to better display the data, but this choice prevents several negative predictions by ordinary cokriging from appearing (4 sites for Zn(DTPA)+OC, 2 sites for Zn(DTPA)+pH and 16 sites for Zn(DTPA)+OC+pH).

(Co)kriging on predictor set data, transformed to reduce positive skewness, produced better estimates of the median concentration of Zn(DTPA) in the testing set than did similar calculations based on untransformed data. For example, the true median of 0.72 mg kg^{-1} was reasonably well predicted by (co)kriging on transformed data (medians of $0.68\text{--}0.74 \text{ mg kg}^{-1}$), but was substantially over-predicted by (co)kriging on untransformed data (medians of $0.82\text{--}0.84 \text{ mg kg}^{-1}$). On the other hand, the overall mean for Zn(DTPA) (1.01 mg kg^{-1}) in the positively skewed testing site data was more successfully reproduced by the more skewed results from (co)kriging on untransformed data (means of $0.93\text{--}0.96 \text{ mg kg}^{-1}$).

The predictions of Zn(DTPA) derived from the kriging models were generally less diverse than the actual values measured at testing sites, in large part because of the averaging process inherent in the kriging and in part because of the effects of back-transformation. This restriction is demonstrated in Table 3 by the compressions in the overall range for predicted values (max–min), the interquartile range ($q_{0.75}\text{--}q_{0.25}$) and the S.D. This compression in the distribution of the predicted values was ameliorated by the incorporation of auxiliary variables. The combination of OC and pH was generally more successful in this regard than either of these secondary variables alone. More detail is provided by Fig. 5, which shows that all methods tended to overestimate low concentrations

(most values are below the 1:1 line) while underestimating high concentrations (most values are above the 1:1 line). The conditional bias demonstrated by this counterclockwise rotation of the elongated cluster of points (Goovaerts, 1997, p. 182) varied substantially among the different kriging approaches. The degree of bias was reduced by the incorporation of auxiliary variables in OCK, OCK_{LG} and OCK_{NS}, while effects for OCK_{RK} were variable.

The quantitative success of prediction was assessed in several ways. The RMSE for predicted concentrations of Zn(DTPA) was decreased when OC or pH was used as a secondary variable. OC was more effective than pH, but OC and pH together almost always provided the lowest RMSE (Table 3). The success of predictions for Zn(DTPA) differed with the type of kriging method used, with the best estimates usually achieved by (co)kriging using transformed data. Differences among OCK_{NS}, OCK_{RK} and OCK_{LG} were small. The trends in calculated r_p^2 for these approaches were similar to those for RMSE (Fig. 5, Table 3).

None of the kriging methods were successful in predicting the highest concentrations of Zn(DTPA) in this skewed data set. About one-half of the total RMSE in predictions was contributed by the two deciles of sites containing the highest concentrations of Zn(DTPA), as shown in Fig. 6 for all four kriging approaches using OC and pH as auxiliary variables. Predictions of Zn(DTPA) by all approaches were more successful for sites in the

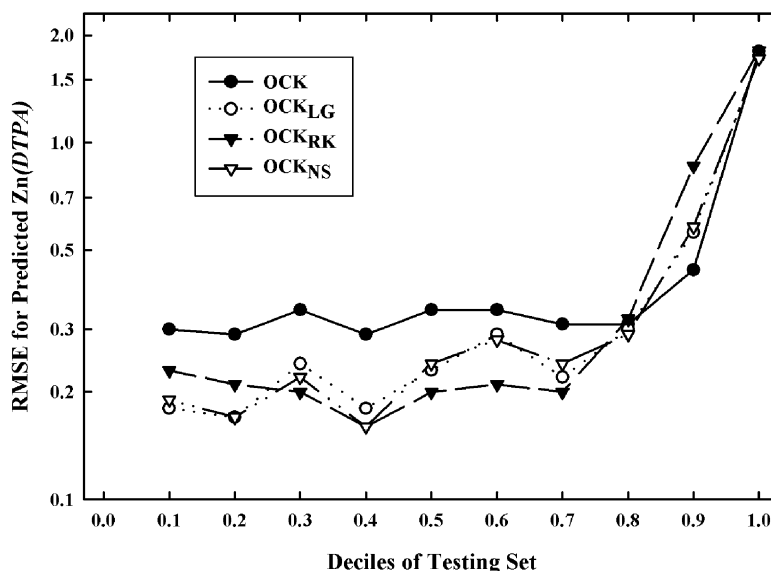


Fig. 6. Root mean squared errors (RMSE) for cokriging predictions of Zn(DTPA), using OC and pH as auxiliary variables. The RMSEs for predictions are shown for population deciles of the 294 testing set sites, ranked by their measured Zn(DTPA) concentrations. Ordinary cokriging (OCK), log-normal ordinary cokriging (OCK_{LG}), rank order ordinary cokriging (OCK_{RK}) and normal score ordinary cokriging (OCK_{NS}) are compared.

lower eight deciles of concentration, and predictions by OCK_{NS} , OCK_{RK} and OCK_{LG} were consistently more accurate than those by OCK .

Making successful predictions of the need of soils for Zn fertilization is one of the goals of understanding the distribution of Zn(DTPA) concentrations. For soil testing purposes, a concentration of 0.5 mg kg^{-1} of Zn(DTPA) is commonly used in North Dakota, and several other States in the mid-western and western USA, as a limit to identify soils which commonly fail to provide adequate Zn for most row crops and small grains (e.g., see numerous on-line recommendations of State Agricultural Extension Services). Predictions by the four kriging approaches were compared for the group of 92 testing set soils, which contained less than 0.5 mg kg^{-1} of Zn(DTPA). Table 3 shows that the RMSEs for OCK_{NS} , OCK_{RK} and OCK_{LG} were commonly one quarter to one-half lower than for OCK , demonstrating that kriging estimates for Zn(DTPA) in these soils, with low plant-available Zn, were most successful when the data for Zn(DTPA) were transformed prior to kriging. The results show also that use of auxiliary variables in kriging was important in reducing RMSE values in these low-Zn soils, especially so when both OC and pH were used together. The MEs also declined in magnitude (although all were negative because measured values were limited by analytical sensitivity while predicted values were not).

As a confirmation of the results presented above, we reversed the testing and predictor sets, and then repeated all kriging inferences. The results were similar (data not shown), leading to the same conclusions with respect to the benefits from cokriging and data transformation.

4. Discussion and conclusions

4.1. Kriging and cokriging with auxiliary information

Ordinary kriging is the anchoring algorithm of geostatistics (Deutsch and Journel, 1998, p. 66). OK is a robust approach and it combines many of the desirable features of alternative methods of geospatial inference (Isaaks and Srivastava, 1989, pp. 318–321). Ordinary cokriging has the advantages of OK, but, in addition, OCK is able to incorporate auxiliary information to further improve estimates of a primary variable (Isaaks and Srivastava, 1989, p. 399). Not surprisingly, the advantages of cokriging with a correlated variable are most evident when the secondary variable has been or will be sampled more densely than the primary variable (Goovaerts, 1999). Examples of cokriging with one

secondary variable can be found easily in the soil literature. For example, we improved estimates of soil Cu in North Dakota soils using auxiliary data for CEC (Wu et al., 2003); McBratney and Webster (1983) interpolated and mapped the silt content of topsoils with the aid of silt or sand content of the corresponding subsoils; Vauclin et al. (1983) made improved estimates of water content from data on the distribution of soil texture; and Stein et al. (1988) predicted the soil moisture deficit using the mean highest water-table depth as an inexpensively measured co-variable.

Although cokriging with more than one secondary variable has the potential to further improve estimates of a primary variable, this approach has been used only occasionally. For example, Han et al. (2003) used clay and silt contents as secondary variables for predictions of N, P and K in soils, but found only slight improvements over OK alone, and only then at their higher sampling density. Using a test data set, Goovaerts (1998) showed that cokriging with two or three auxiliary variables produced better estimates than kriging with none, but did not report whether the use of more than one auxiliary variable increased benefits above use of a single variable.

Our results for Zn(DTPA) demonstrated a clear benefit from cokriging with two auxiliary variables in predicting Zn(DTPA). This improvement was enhanced by several factors. First, by assessing this benefit at testing set locations where the auxiliary variables were available, the benefit was maximized and most easily demonstrated, as was our intent. Improvements at the more numerous sites having no auxiliary data would clearly be less. Secondly, the benefit from using two auxiliary variables was enhanced by the fact that OC and pH were themselves not highly correlated, even though each was correlated with Zn(DTPA). Thirdly, benefits from use of auxiliary variables were favored by our comparison of predictor and testing populations of equivalent size. And, finally, the spatial dependency for Zn(DTPA) in our data was not particularly strong, so that there was ample opportunity for improvement from auxiliary information.

In a chemical sense, the benefits from using pH in estimating the level of available Zn in these soils arose, presumably, from the well-recognized influence of pH on metal solubility and extractability. Similarly, benefits from using OC as an auxiliary variable arose, presumably, from the contributions to extractable Zn(DTPA) from Zn associated with the metal-binding sites on soil organic matter. In addition to providing mechanistically reasonable explanations for their association with extractable soil Zn, a separate advantage to using pH and OC as auxiliary variables is that data are often readily

available from other sources, such as state soil testing programs or the compiled characterization data for soil pedons (National Soil Survey Center, 2002).

The main disadvantage to cokriging is that it is more computationally demanding than kriging. If n variables are involved, $n \times (n+1)/2$ auto- and cross-semivariograms must be inferred and jointly modeled, and a large cokriging system must be solved (Goulard and Voltz, 1992; Goovaerts, 1999). Thus, in our work with Zn(DTPA), OC and pH, we needed to model jointly six variograms; and, to reduce computational time, the cokriging calculations were carried out on a high speed computing system at the Cornell Theory Center (Cornell University, Ithaca, NY).

4.2. Kriging with transformed data

Ordinary (co)kriging is the best linear unbiased estimator of variables at unmeasured locations under the stationarity assumption (Goovaerts, 1997). However, real data hardly ever completely satisfy this assumption, as illustrated by our data for Zn(DTPA), which were highly skewed. The resulting experimental variogram for Zn(DTPA) was erratic, making it hard to see the spatial structure or fit a model. For these data, logarithmic, rank order and normal score transformations were successful in removing most skewness, and the recomputed variograms exhibited clear spatial structure.

While transformation of skewed data may improve its suitability for (co)kriging, transformation and back-transformation may have other effects that are hard to interpret or may add uncertainty. As already mentioned, all three of the back-transformation methods used in this paper yield a prediction, in the original units of measure, of a median rather than a mean value for each (co)kriged location. In addition, the log transformation is reported to be especially sensitive to slight variations in the sill of the variogram and exponentiation during back-transformation exaggerates variability (Armstrong and Boufassa, 1988; Deutsch and Journel, 1998, pp. 75–76; Roth, 1998). Despite these concerns, OK_{LG} and OCK_{LG} were among the most successful methods for predicting the concentrations of Zn(DTPA), as shown by Table 3 and Fig. 5.

The rank order and normal score transformations share the advantage that data of diverse types, scales, ranges and reliability can be integrated in the kriging process (Journel and Deutsch, 1997). The normal score transformation has the added advantage of producing a normalized distribution out of any rank-ordered data, but the inherent assumption that the distribution is truly multi-Gaussian cannot be adequately tested (Goovaerts,

1999). The rank order approach has the added limitation of generating only a uniform distribution of transformed observations, and it has the potential to produce some estimates that must be adjusted arbitrarily to fall within the permissible standardized range for ranks [0,1] (Journel and Deutsch, 1997), a problem illustrated by our own results. Nevertheless, both rank order and normal score transformation of the Zn(DTPA) data of this study yielded variograms with clear spatially dependent structure as shown in Fig. 4, and OCK_{RK} and OCK_{NS} were both better predictors for Zn(DTPA) than OCK .

In summary, it is clear that predictions of Zn(DTPA) in soils at sparsely distributed sites in northern North Dakota were improved substantially by cokriging with auxiliary data for OC or pH. Utilizing OC and pH together as secondary variables improved the predictions further. We believe that predictions of crop-available forms of other trace metals in soils are likely to benefit similarly, because OC is related to the abundance of metal binding sites and pH is well recognized as influencing metal solubility. Transformation of our Zn(DTPA) data to remove skewness and reduce the distorting influence of high outliers improved estimates of Zn(DTPA) at testing locations where these data were treated as unavailable during kriging or cokriging. Logarithmic, rank order and normal score transformations were all about equally successful in this regard. Predictions of Zn(DTPA) in soils containing low levels of extractable Zn were especially improved by transformation prior to cokriging.

Acknowledgments

We are very grateful for the valuable comments and suggestions provided by two anonymous reviewers and Editor A.B. McBratney. We also thank Dr. A.G. Journel at Stanford University and Dr. P. Goovaerts at the University of Michigan for their advice. This research was conducted partially using the resources of the Cornell Theory Center, which received funding from Cornell University, New York State, federal agencies, foundations and corporate partners.

References

- Armstrong, M., Boufassa, A., 1988. Comparing the robustness of ordinary kriging and lognormal kriging—outlier resistance. *Math. Geol.* 20, 447–457.
- Cattle, J.A., McBratney, A.B., Minasny, B., 2002. Kriging method evaluation for assessing the spatial distribution of urban soil lead contamination. *J. Environ. Qual.* 31, 1576–1588.
- Deutsch, C.V., Journel, A.G., 1998. *GSLIB, Geostatistical Software Library and User's Guide*. Oxford University Press, New York, NY, p. 369.

- Franzen, D.W., 1999. North Dakota survey of soil copper, pH, zinc, and boron levels. In: NDSU Ext. Service (Ed.), Extension Report, vol. 52. Ext. Service, North Dakota State Univ., Fargo, ND.
- Gibson, R.S., 1994. Zinc nutrition in developing countries. *Nutr. Res. Rev.* 7, 151–173.
- Goovaerts, P., 1997. *Geostatistics for Natural Resources Evaluation*. Oxford University Press, New York, NY, pp. 483.
- Goovaerts, P., 1998. Ordinary cokriging revisited. *Math. Geol.* 30, 21–42.
- Goovaerts, P., 1999. Geostatistics in soil science: state-of-the-art and perspectives. *Geoderma* 89, 1–45.
- Goulard, M., Voltz, M., 1992. Linear coregionalization model: tools for estimation and choice of cross-variogram matrix. *Math. Geol.* 24, 269–286.
- Han, S., Schneider, S.M., Evans, R.G., 2003. Evaluating cokriging for improving soil nutrient sampling efficiency. *Trans. Am. Soc. Agric. Eng.* 46, 845–849.
- Holmgren, G.G.S., Meyer, M.W., Chaney, R.L., Daniels, R.B., 1993. Cadmium, lead, zinc, copper, and nickel in agricultural soils of the United States of America. *J. Environ. Qual.* 22, 335–348.
- Isaaks, E.H., Srivastava, R.M., 1989. *Applied Geostatistics*. Oxford University Press, New York, NY, pp. 561.
- Journel, A.G., 1980. The lognormal approach to predicting local distributions of selective mining unit grades. *Math. Geol.* 12, 285–303.
- Journel, A.G., Deutsch, C.V., 1997. Rank order geostatistics: a proposal for a unique coding and common processing of diverse data. In: Schofield, E.Y.B.a.N.A. (Ed.), *Geostatistics Wollongong '96*. Kluwer Academic Publishers, Dordrecht, The Netherlands, pp. 174–187.
- Juang, K.W., Lee, D.Y., 1998. A comparison of three kriging methods using auxiliary variables in heavy-metal contaminated soils. *J. Environ. Qual.* 27, 355–363.
- Juang, K.W., Lee, D.Y., Ellsworth, T.R., 2001. Using rank-order geostatistics for spatial interpolation of highly skewed data in a heavy-metal contaminated site. *J. Environ. Qual.* 30, 894–903.
- Lindsay, W.L., Norvell, W.A., 1978. Development of a DTPA test for zinc, iron, manganese, and copper. *Soil Sci. Soc. Am. J.* 42, 421–428.
- McBratney, A.B., Webster, R., 1983. Optimal interpolation and isarithmic mapping of soil properties: V. Coregionalization and multiple sampling strategy. *J. Soil Sci.* 34, 137–162.
- McBratney, A.B., Webster, R., McLaren, R.G., Spiers, R.B., 1982. Regional variation of extractable copper and cobalt in the topsoil of south-east Scotland. *Agronomie* 2, 969–982.
- National Soil Survey Center, 2002. National soil characterization data base. Soil Survey Division, USDA-NRCS.
- Norvell, W.A., Wu, J., Hopkins, D.G., Welch, R.M., 2000. Association of cadmium in durum wheat grain with soil chloride and chelate-extractable soil cadmium. *Soil Sci. Soc. Am. J.* 65, 2162–2168.
- Petersen, R.G., Calvin, L.D., 1986. Sampling. In: A. Klute (Editor), *Methods of soil analysis: Part 1. Agron. Monogr.* 9. ASA and SSSA, Madison, WI, pp. 33–51.
- Roth, C., 1998. Is lognormal kriging suitable for local estimation? *Math. Geol.* 30, 999–1009.
- Saito, H., Goovaerts, P., 2000. Geostatistical interpolation of positively skewed and censored data in a dioxin-contaminated site. *Environ. Sci. Technol.* 34, 4228–4235.
- Sillanpää, M., 1990. Micronutrient assessment at the country level: an international study. *FAO Soils Bulletin*, vol. 63. FAO/Finnish International Development Agency, Rome, Italy.
- Stein, A., Hoogererf, M., Bouma, J., 1988. Use of soil-map delineations to improve (co)kriging of point data on moisture deficits. *Geoderma* 43, 163–177.
- Van Meirvenne, M., Goovaerts, P., 2001. Evaluating the probability of exceeding a site-specific soil cadmium contamination threshold. *Geoderma* 102, 75–100.
- Vauclin, M., Vieira, S.R., Vachaud, G., Nielsen, D.R., 1983. The use of cokriging with limited field soil observations. *Soil Sci. Soc. Am. J.* 47, 175–184.
- Webster, R., Oliver, M.A., 2001. *Geostatistics for Environmental Scientists*. John Wiley & Sons Inc., New York, NY. 271 pp.
- Welch, R.M., 1995. Micronutrient nutrition of plants. *Crit. Rev. Plant Sci.* 14, 49–82.
- Wu, J., Norvell, W.A., Hopkins, D.G., Smith, D.B., Ulmer, M.G., Welch, R.M., 2003. Improved prediction and mapping of soil copper by kriging with auxiliary data for CEC. *Soil Sci. Soc. Am. J.* 67, 919–927.